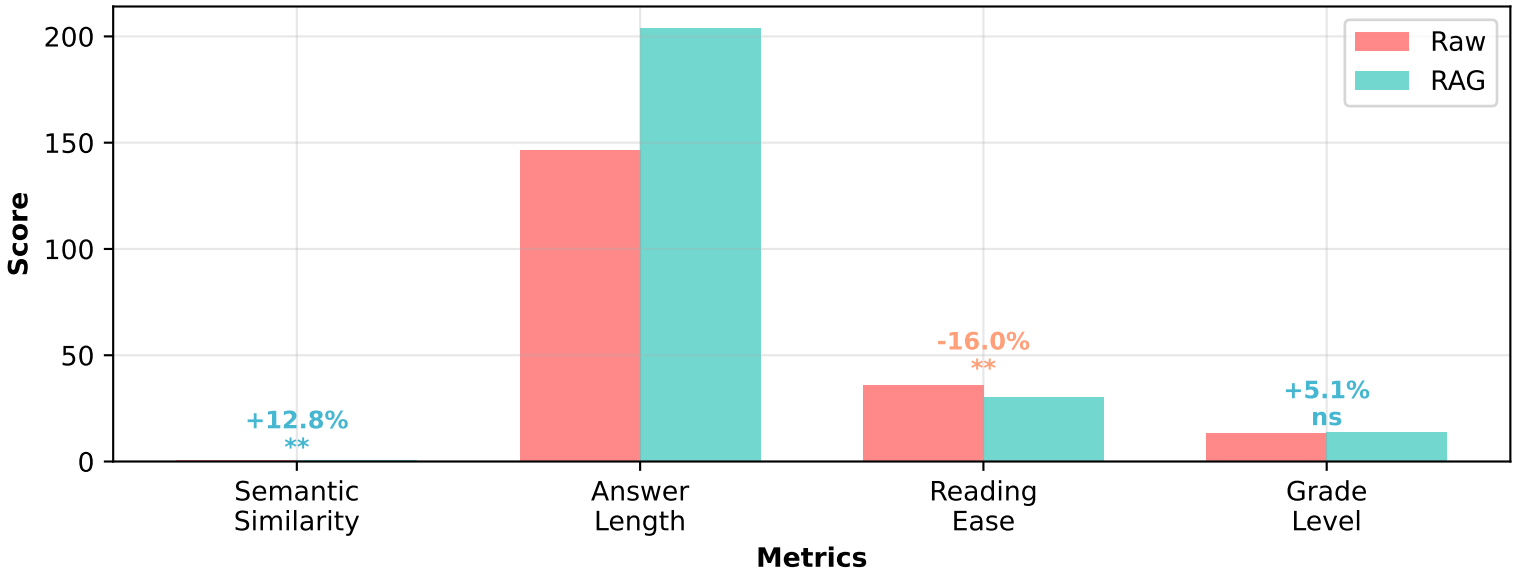


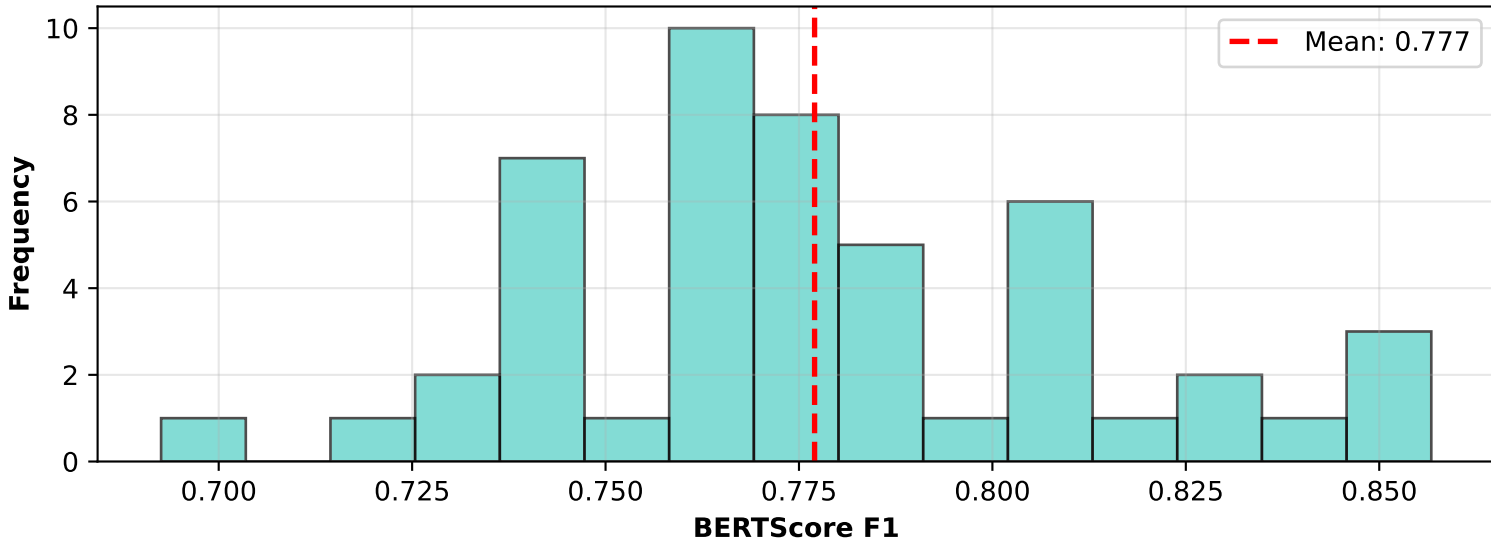
# CGT Chatbot Performance Analysis: RAG vs Raw Approach

## Comprehensive Results Dashboard

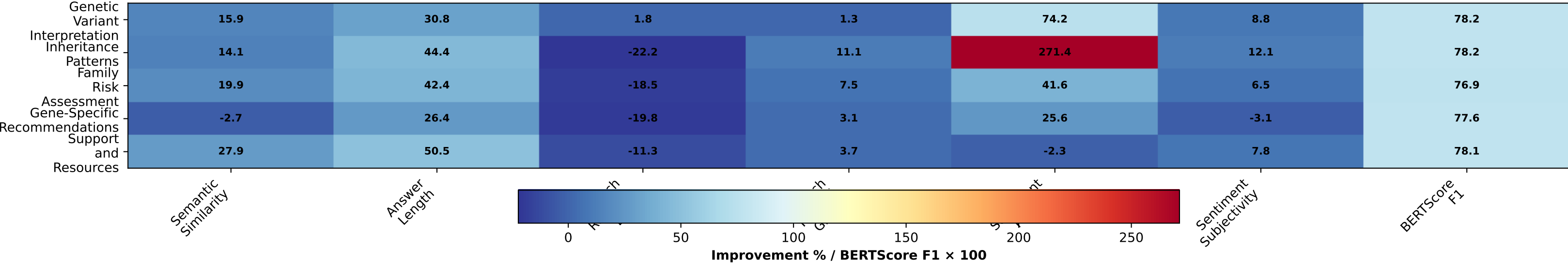
Overall Performance Comparison



BERTScore F1 Distribution (RAG vs Raw Semantic Similarity)



Performance by Question Category (% Improvement / BERTScore F1 × 100)



Overall Performance Statistics

Category Performance Summary

Metric	Raw Mean	RAG Mean	Improvement %	P-value	Significant
Semantic Similarity	0.521	0.588	+12.8%	0.0043	Yes
Answer Length	146.735	203.878	+38.9%	0.0000	Yes
Flesch Reading Ease	35.898	30.161	-16.0%	0.0050	Yes
Flesch Kincaid Grade	13.376	14.053	+5.1%	0.0722	No
BERTScore F1	N/A	0.777	N/A	N/A	49/49

Question Category	N Questions	Semantic Sim %	BERTScore F1	Significant Metrics
Genetic Variant Interpretation	5	+15.9%	0.782	1/6
Inheritance Patterns	7	+14.1%	0.782	2/6
Family Risk Assessment	11	+19.9%	0.769	2/6
Gene-Specific Recommendations	14	-2.7%	0.776	1/6
Support and Resources	12	+27.9%	0.781	2/6